

Huawei Cloud EulerOS

常见问题

文档版本 01
发布日期 2024-06-30



版权所有 © 华为云计算技术有限公司 2024。保留一切权利。

非经本公司书面许可，任何单位和个人不得擅自摘抄、复制本文档内容的部分或全部，并不得以任何形式传播。

商标声明



HUAWEI和其他华为商标均为华为技术有限公司的商标。

本文档提及的其他所有商标或注册商标，由各自的所有人拥有。

注意

您购买的产品、服务或特性等应受华为云计算技术有限公司商业合同和条款的约束，本文档中描述的全部或部分产品、服务或特性可能不在您的购买或使用范围之内。除非合同另有约定，华为云计算技术有限公司对本文档内容不做任何明示或暗示的声明或保证。

由于产品版本升级或其他原因，本文档内容会不定期进行更新。除非另有约定，本文档仅作为使用指导，本文档中的所有陈述、信息和建议不构成任何明示或暗示的担保。

目录

1 CentOS Linux 停止维护后如何应对?	1
2 华为云针对 CentOS EOL 有没有迁移方案?	3
3 如何安装 mlnx 驱动?	5
4 如何开启 HCE OS 操作系统的 SELinux 功能?	8
5 迁移系统后, 如何更改控制台操作系统名称?	10
6 Huawei Cloud EulerOS、openEuler 和 EulerOS 镜像的主要区别是什么?	13
7 如何打开内核 wireguard 模块以及安装 wireguard-tools?	14
8 如何将 docker 工具的用户凭证保存方式配置成与社区一致?	16
9 OOM 相关参数配置与原因排查.....	17
10 IPVS 报错问题说明.....	22
11 ECS 开启 IPv6 后, HCE OS 系统内无法获取到 IPv6 地址.....	24

1 CentOS Linux 停止维护后如何应对?

CentOS官方已计划停止维护CentOS操作系统，华为云上CentOS公共镜像来源于CentOS官方，当CentOS操作系统停止维护后，华为云将会同时停止对该操作系统的支持。本文主要介绍CentOS操作系统停止维护带来的影响，并针对影响提供应对策略。

背景信息

2020年12月08日，CentOS官方宣布了停止维护CentOS Linux的计划，并推出了CentOS Stream项目。更多信息，请参见[CentOS官方公告](#)。

CentOS 8系统2021年12月31日已停止维护服务，CentOS 7系统将于2024年06月30日停止维护服务。CentOS官方不再提供CentOS 9及后续版本，不再支持新的软件和补丁更新。CentOS用户现有业务随时面临宕机和安全风险，并无法确保及时恢复。

影响

基于CentOS官方的变更计划，对CentOS操作系统的使用者产生的影响如下所述：

- 2021年12月31日以后，CentOS 8的使用者将无法获得包括问题修复和功能更新在内的任何软件维护和支持。
- 2024年06月30日以后，CentOS 7的使用者将无法获得包括问题修复和功能更新在内的任何软件维护和支持。

对于华为云的公共镜像及服务支持存在一定影响：

- 华为云暂不会下线CentOS 8公共镜像，同时已经使用CentOS 8创建的ECS实例运行不会受到影响，但将停止更新镜像。
- 华为云对于CentOS操作系统的服务支持将和CentOS官方日期保持同步。2021年12月31日以后将不再对CentOS 8提供服务支持；对CentOS 7的服务支持将持续至2024年6月30日。

应对策略

为了保障使用CentOS系统的业务正常运行，华为云为您提供替换CentOS操作系统的应对策略。替换CentOS操作系统的方式分为两类，切换操作系统和迁移操作系统。

切换到Debian或Ubuntu具体操作、详见[切换操作系统](#)。

- 将CentOS操作系统切换为[支持切换的操作系统](#)。

如果现有的ECS配置（网卡、磁盘、VPN等配置的类型和数量）都不需要改变，仅需要修改ECS的操作系统镜像，并且您的软件和原操作系统耦合度较低，建议使用系统切换。

- 切换到Huawei Cloud EulerOS具体操作，详见[将操作系统切换为HCE OS](#)。
- 切换到Debian或Ubuntu具体操作、详见[切换操作系统](#)。
- 将CentOS操作系统迁移为Huawei Cloud EulerOS操作系统。
如果现有的ECS配置（网卡、磁盘、VPN等配置的类型和数量）都不需要改变，希望保留操作系统软件的配置参数，可以通过操作系统迁移的方式迁移到Huawei Cloud EulerOS。
系统迁移详见[将操作系统迁移为HCE OS](#)。

系统切换和迁移的区别如下表，请根据需要选择合适的替换方式。

表 1-1 系统切换和迁移的区别

区别	系统切换	系统迁移
数据备份	<ul style="list-style-type: none"> ● 切换操作系统会清除系统盘数据，包括系统盘上的系统分区和所有其它分区。 ● 切换操作系统不影响数据盘数据。 	<ul style="list-style-type: none"> ● 迁移操作系统不会清除系统盘数据，为避免系统软件的数据丢失，建议将其备份。 ● 迁移操作系统不影响数据盘数据。
个性化设置	切换操作系统后，当前操作系统内的个性化设置（如DNS、主机名等）将被重置，需重新配置。	迁移操作系统后，当前操作系统内的个性化设置（如DNS、主机名等）不需重新配置。

表 1-2 支持切换的操作系统

操作系统	概述	适用人群
Huawei Cloud EulerOS	<p>Huawei Cloud EulerOS（简称HCE OS）是基于openEuler构建的云上操作系统。</p> <p>HCE OS打造云原生、高性能、高安全、易迁移等能力，加速用户业务上云，提升用户的应用创新空间，可替代CentOS、EulerOS等公共镜像。</p>	适用于希望使用免费镜像，并延续开源社区镜像使用习惯的个人或企业。
Debian、Ubuntu操作系统	Linux的其他发行版操作系统，不同操作系统在使用习惯和应用兼容性上存在一定差异。	适用于可以自行应对操作系统切换成本的个人或企业。

2 华为云针对 CentOS EOL 有没有迁移方案?

背景信息

CentOS 8系统2021年12月31日已停止维护服务，CentOS 7系统将于2024年06月30日停止维护服务。CentOS不再支持新的软件和补丁更新。CentOS用户现有业务随时面临宕机和安全风险，并无法确保及时恢复。

HCE OS操作系统从云原生混部竞争力、安全可信、快速迁移、高效运维、专业认证等方面为用户提供专业云服务、解决CentOS停服带来的影响。HCE OS提供了迁移工具，可将CentOS、EulerOS等操作系统平滑迁移至HCE OS操作系统。

兼容性评估

华为HCE OS操作系统已具备完整代替CentOS的技术能力，完全自主可控，并基于openEuler开源社区持续自主演进。南向支持6大类400种板卡，基本覆盖主流计算产品。北向100%兼容主流的应用场景（云原生、存储、数据库、大数据、WEB等）。超过5000种应用软件通过兼容性认证，基本能够替代CentOS的各种部署。

为满足CentOS系列到HCE OS搬迁的准确性和安全性，请您使用兼容性工具对待迁移软件快速进行扫描，获取评估结果。

对于可兼容的应用软件，迁移过程中并不会修改软件配置，迁移完成后无需重新配置；对于部分不兼容的应用软件，评估报告给出相应的规避策略，请在迁移之后进行相应的适配。

迁移能力评估

HCE OS已有成熟的搬迁指导，按照分布式集群应用、主备应用、单机应用三种类型对各种应用进行归类，并制定相应的搬迁方案：

- 针对分布式集群软件，如大数据、分布式存储。CentOS搬迁无需中断业务，基于分布式软件伸缩扩容机制，HCE OS实现滚动代替，平滑搬迁。
- 针对主备应用，如数据库等。CentOS搬迁无需中断业务，先备后主，基于主备状态同步机制，平滑搬迁。
- 针对单机应用，CentOS搬迁需中断业务，割接式搬迁，该类搬迁方案和原应用重新部署方式等同，属于成熟方案。

联系我们

华为云专业的服务工程师团队为您提供7*24小时保驾护航服务。若您在使用华为云时遇到问题，您可以通过提交[工单](#)寻求帮助，我们将竭诚为您服务。

3 如何安装 mlnx 驱动?

本节介绍如何在HCE OS 2.0操作系统（包含x86架构和arm架构系统）安装mlnx驱动。

约束与限制

- HCE OS 2.0的内核为linux 5.10及以上版本。
- CX6网卡驱动为23.10-1.1.9.0-LTS及以上版本。

前提条件

已经安装了5.10或更高内核版本的HCE OS 2.0系统。

在 x86 架构安装 mlnx 驱动

1. 下载CX6网卡驱动安装包[MLNX_OFED_LINUX-23.10-1.1.9.0-openeuler22.03-x86_64.tgz](#)。
2. 解压驱动安装包并进入工作目录。

```
tar -xf MLNX_OFED_LINUX-23.10-1.1.9.0-openeuler22.03-x86_64.tgz
cd MLNX_OFED_LINUX-23.10-1.1.9.0-openeuler22.03-x86_64
```

3. 安装CX6网卡驱动软件。

```
./mlnxofedinstall --basic --without-depcheck --distro OPENEULER22.03 --
force --kernel 5.10.0-60.18.0.50.oe2203.x86_64 --kernel-sources /lib/
modules/$(uname -r)/build
```

📖 说明

其中，“5.10.0-60.18.0.50.oe2203.x86_64”是官方MLNX_OFED包本身编译时的内核版本。

4. 创建链接。

```
ln -s /lib/modules/5.10.0-60.18.0.50.oe2203.x86_64/extra/mlnx-
ofa_kernel /lib/modules/$(uname -r)/weak-updates/
```

```
ln -s /lib/modules/5.10.0-60.18.0.50.oe2203.x86_64/extra/kernel-mft /lib/
modules/$(uname -r)/weak-updates/
```

```
depmod -a
```

5. 执行reboot命令重新系统。

6. 执行`/etc/init.d/openibd status`命令查看驱动安装结果。
显示如下信息表示驱动安装成功。

```
[root@lnmp ~]# /etc/init.d/openibd status
HCA driver loaded

The following OFED modules are loaded:

rdma_ucm
rdma_cm
ib_ipoib
mlx5_core
mlx5_ib
ib_uverbs
ib_umad
ib_cm
ib_core
mlxfw

[root@lnmp ~]# █
```

在 arm 架构安装 mlnx 驱动

1. 下载CX6网卡驱动安装包[MLNX_OFED_LINUX-23.10-1.1.9.0-openeuler22.03-aarch64.tgz](#)。
2. 解压驱动安装包并进入工作目录。
`tar -xf MLNX_OFED_LINUX-23.10-1.1.9.0-openeuler22.03-aarch64.tgz`
`cd MLNX_OFED_LINUX-23.10-1.1.9.0-openeuler22.03-aarch64`
3. 安装CX6网卡驱动软件。

```
./mlnxofedinstall --basic --without-depcheck --distro OPENEULER22.03 --force --kernel 5.10.0-60.18.0.50.oe2203.aarch64 --kernel-sources /lib/modules/$(uname -r)/build
```

📖 说明

其中，“5.10.0-60.18.0.50.oe2203.aarch64”是官方MLNX_OFED包本身编译时的内核版本。

4. 执行如下命令创建链接。
`ln -s /lib/modules/5.10.0-60.18.0.50.oe2203.aarch64/extra/mlnx-ofa_kernel /lib/modules/$(uname -r)/weak-updates/`
`ln -s /lib/modules/5.10.0-60.18.0.50.oe2203.aarch64/extra/kernel-mft /lib/modules/$(uname -r)/weak-updates/`
`depmod -a`
5. 执行`reboot`命令重新系统。
6. 执行`/etc/init.d/openibd status`命令查看驱动安装结果。
显示如下信息表示驱动安装成功。

```
[root@lnmp ~]# /etc/init.d/openibd status  
  
HCA driver loaded  
  
The following OFED modules are loaded:  
  
rdma_ucm  
rdma_cm  
ib_ipoib  
mlx5_core  
mlx5_ib  
ib_uverbs  
ib_umad  
ib_cm  
ib_core  
mlxfw  
  
[root@lnmp ~]# █
```

4 如何开启 HCE OS 操作系统的 SELinux 功能?

HCE OS操作系统默认关闭SELinux功能。如果业务需要开启SELinux 功能，请参照本节指导操作。

⚠ 注意

请按照本节指导开启SELinux功能，勿直接通过/etc/selinux/config开启SELinux功能，否则可能会出现无法登录的问题。

操作步骤

1. 打开配置文件/boot/grub2/grub.cfg，删除selinux=0。
2. 执行touch /.autorelabel命令。
/.autorelabel文件将触发OS在启动过程中对磁盘上所有文件relabel重新打selinux标签，该过程可能需要持续几分钟。relabel完成后OS将自动重启一次并生效，同时自动删除/.autorelabel文件确保下次不会再重复执行relabel动作。
3. 打开配置文件/etc/selinux/config，设置SELINUX=permissive，并执行reboot重启操作系统。

```
# This file controls the state of SELinux on the system.
# SELINUX= can take one of these three values:
#   enforcing - SELinux security policy is enforced.
#   permissive - SELinux prints warnings instead of enforcing.
#   disabled - No SELinux policy is loaded.
SELINUX=permissive
# SELINUXTYPE= can take one of three values:
#   targeted - Targeted processes are protected,
#   minimum - Modification of targeted policy. Only selected processes are protected.
#   mls - Multi Level Security protection.
SELINUXTYPE=targeted
```

4. 再次打开配置文件/etc/selinux/config，设置SELINUX=enforcing，并执行reboot重启操作系统。

```
# This file controls the state of SELinux on the system.
# SELINUX= can take one of these three values:
#   enforcing - SELinux security policy is enforced.
#   permissive - SELinux prints warnings instead of enforcing.
#   disabled - No SELinux policy is loaded.
SELINUX=enforcing
# SELINUXTYPE= can take one of three values:
#   targeted - Targeted processes are protected,
#   minimum - Modification of targeted policy. Only selected processes are protected.
#   mls - Multi Level Security protection.
SELINUXTYPE=targeted
```

5. 重启后执行**getenforce**命令查看SELinux状态。
显示Enforcing表示SELinux已经开启。

```
[root@localhost ~]#
[root@localhost ~]# getenforce
Enforcing
```

5 迁移系统后，如何更改控制台操作系统名称？

问题背景

原操作系统（例如CentOS 7.9）迁移为Huawei Cloud EulerOS后，控制台仍然显示原操作系统名称CentOS 7.9而不是Huawei Cloud EulerOS。

您可通过创建私有镜像、再切换到此私有镜像的方式，将控制台操作系统名称更改为Huawei Cloud EulerOS。



操作步骤

1. 登录[ECS控制台](#)。
2. 在待迁移系统的弹性云服务器的“操作”列下，选择“更多 > 镜像 > 创建镜像”。
3. 在“创建私有镜像”页面，配置如下镜像信息。
 - 区域：服务器所在区域，请保持此默认配置。
 - 创建方式：创建私有镜像，请保持此默认配置。
 - 镜像类型：系统盘镜像，请保持此默认配置。
 - 镜像源：迁移系统的弹性云服务器，请保持此默认配置。
 - 名称：填写便于识别的镜像名称。
 - 协议：阅读并勾选协议。

目前镜像服务已进入商业化阶段，私有镜像会收取一定的存储费用。详细计费标准可参考镜像服务计费标准

镜像类型和来源

* 区域: 华东-上海一
不同区域的云服务器产品之间内网互不相通；请就近选择靠近您业务的区域，可减少网络时延，提高访问速度。

* 创建方式: **创建私有镜像** 导入私有镜像

* 镜像类型: **系统盘镜像** 整机镜像 数据盘镜像

* 选择镜像源: **云服务器** 裸金属服务器

当前关机或开机状态的弹性云服务器可以用来创建私有镜像。
创建镜像前，请确保弹性云服务器已完成相关配置。了解更多
请勿在创建镜像过程中对所选择的弹性云服务器及其关联资源进行其他操作。

所有状态 ID fc62clee-7f39-49de-985 X Q C

名称	操作系统	运行状态	私有IP地址	创建时间
Servers-to-be-migrated	CentOS 7.9 64bit	运行中		2023/09/14 11:16:25 GM...

当前选择: Servers-to-be-migrated | 操作系统: CentOS 7.9 64bit | 系统盘: 通用型SSD | 40 GiB
[购买弹性云服务器](#)

配置信息

加密: 未加密

* 名称:

* 企业项目: --请选择-- C

标签: 如果您需要使用同一标签标识多种云资源，即所有服务均可在标签输入框下拉选择同一标签。建议在TMS中创建预定义标签。[查看预定义标签](#)
标签键: 标签值:
您还可以添加10个标签。

描述:
0/1,024

协议: 我已阅读并同意《[镜像制作承诺书](#)》和《[镜像免责声明](#)》

立即创建

4. 单击“立即创建”创建私有镜像。
5. 确认镜像信息，单击“提交”。
6. 返回弹性云服务器控制台，在待切换操作系统的弹性云服务器的“操作”列下，选择“更多 > 镜像 > 切换操作系统”。
7. 在“切换操作系统”界面配置如下参数。
 - 勾选立即关机。
 - 镜像：选择私有镜像。
 - 登录凭证：选择“使用镜像密码”。



- 单击“确定”并根据界面提示完成验证。
- 阅读并勾选声明，单击“确定”。
系统切换后，控制台系统名称即更改为Huawei Cloud EulerOS。

6 Huawei Cloud EulerOS、openEuler 和 EulerOS 镜像的主要区别是什么？

Huawei Cloud EulerOS、openEuler和EulerOS镜像均为华为自研镜像，主要区别如下表 1所示：

表 6-1 Huawei Cloud EulerOS、openEuler 和 EulerOS 镜像的区别与联系

镜像类型	描述
Huawei Cloud EulerOS (简称HCE OS)	Huawei Cloud EulerOS (简称HCE OS)，是基于openEuler开发的一款商业发行版镜像，可替代CentOS、EulerOS等操作系统，并提供专业的维护保障能力，镜像目前免费对用户使用。 说明 Huawei Cloud EulerOS 2.0是基于openEuler 22.03 LTS版本构建的云上操作系统。
openEuler	openEuler是一款开源镜像，您可以免费使用，但是不提供商业维护保障能力。openEuler最初由华为研发，但是在2021年11月9日正式捐赠给开放原子开源基金会，openEuler的技术支持由开源社区提供。
EulerOS	EulerOS是基于开源技术的企业级Linux操作系统软件，具备高安全性、高可扩展性、高性能等技术特性，能够满足客户IT基础设施和云计算服务等多业务场景需求。 说明 EulerOS是基于开源操作系统openEuler进行开发的华为内部的操作系统。

7 如何打开内核 wireguard 模块以及安装 wireguard-tools?

说明

wireguard-tools工具来源于社区，如果您在使用中遇到问题，可通过<https://github.com/WireGuard/wireguard-tools/pulls>获取帮助。

打开内核 wireguard 模块

您通过命令 `modprobe wireguard` 打开内核 wireguard 模块。

安装 wireguard-tools

步骤1 执行以下命令安装依赖。

```
dnf install gcc make
```

步骤2 执行以下命令下载 wireguard-tools 源码包。

```
wget https://git.zx2c4.com/wireguard-tools/snapshot/wireguard-  
tools-1.0.20210914.tar.xz
```

步骤3 执行以下命令解压上述源码包。

```
tar -xf wireguard-tools-1.0.20210914.tar.xz
```

步骤4 进入 wireguard-tools-1.0.20210914/src 目录，依次执行以下命令编译安装。

```
make
```

```
make install
```

步骤5 验证安装是否成功。

可以执行 `wg -h` 和 `wg-quick -h` 命令验证是否安装成功，如图所示。

```
[root@localhost ~]# wg -h
Usage: wg <cmd> [<args>]

Available subcommands:
  show: Shows the current configuration and device information
  showconf: Shows the current configuration of a given WireGuard interface, for use with `setconf`
  set: Change the current configuration, add peers, remove peers, or change peers
  setconf: Applies a configuration file to a WireGuard interface
  addconf: Appends a configuration file to a WireGuard interface
  syncconf: Synchronizes a configuration file to a WireGuard interface
  genkey: Generates a new private key and writes it to stdout
  genspk: Generates a new preshared key and writes it to stdout
  pubkey: Reads a private key from stdin and writes a public key to stdout
You may pass `--help' to any of these subcommands to view usage.
[root@localhost ~]# wg-quick -h
Usage: wg-quick [ up | down | save | strip ] [ CONFIG_FILE | INTERFACE ]

CONFIG_FILE is a configuration file, whose filename is the interface name
followed by `.conf'. Otherwise, INTERFACE is an interface name, with
configuration found at /etc/wireguard/INTERFACE.conf. It is to be readable
by wg(8)'s `setconf' sub-command, with the exception of the following additions
to the [Interface] section, which are handled by wg-quick:

- Address: may be specified one or more times and contains one or more
  IP addresses (with an optional CIDR mask) to be set for the interface.
- DNS: an optional DNS server to use while the device is up.
- MTU: an optional MTU for the interface; if unspecified, auto-calculated.
- Table: an optional routing table to which routes will be added; if
  unspecified or `auto', the default table is used. If `off', no routes
  are added.
- PreUp, PostUp, PreDown, PostDown: script snippets which will be executed
  by bash(1) at the corresponding phases of the link, most commonly used
  to configure DNS. The string `%i' is expanded to INTERFACE.
- SaveConfig: if set to `true', the configuration is saved from the current
  state of the interface upon shutdown.

See wg-quick(8) for more info and examples.
[root@localhost ~]#
```

---结束

8 如何将 docker 工具的用户凭证保存方式配置成与社区一致?

问题背景

社区版本的docker工具，使用docker login命令登录成功后，会将用户的用户名、密码等数据以base64的格式保存在用户配置文件，存在较大安全隐患，所以HCE OS 2.0提供的docker工具，将默认的保存方式改为了加密保存。部分社区工具暂时不支持该安全特性，需要手动将保存方式改为社区的保存方案。

如何将凭证保存方式修改为社区方案

1. 配置环境变量
export USE_DECRYPT_AUTH=true
2. 使用docker login命令重新登录
docker login

```
[root@ecs-7eca ~]# docker login [REDACTED]
WARNING: Error loading config file: /root/.docker/config.json: Invalid auth configuration file
Username: [REDACTED]
Password:
WARNING! Your password will be stored unencrypted in /root/.docker/config.json.
Configure a credential helper to remove this warning. See
https://docs.docker.com/engine/reference/commandline/login/#credentials-store
Login Succeeded
```

3. 验证完成后，建议将环境变量配置保存在持久文件中（如 ~/.bash_profile，/etc/profile等），以便重启后生效
echo "export USE_DECRYPT_AUTH=true" >> ~/.bash_profile

9 OOM 相关参数配置与原因排查

OOM 相关概念

OOM (Out Of Memory, 简称OOM) 指系统内存已用完, 在linux系统中, 如果内存用完会导致系统无法正常工作, 触发系统panic或者OOM killer。

OOM killer是linux内核的一个机制, 该机制会监控那些占用内存过大的进程, 尤其是短时间内消耗大量内存的进程, 在系统的内存即将不够用时结束这些进程从而保障系统的整体可用性。

OOM 相关参数

表 9-1 OOM 相关参数

参数名称	参数说明	取值	修改方式
panic_on_oom	<p>panic_on_oom参数是控制系统遇到OOM时如何反应的。当系统遇到OOM的时候, 通常会有两种选择:</p> <ul style="list-style-type: none">• 触发系统panic, 可能会出现频繁宕机的情况。• 选择一个或者几个进程, 触发OOM killer, 结束选中的进程, 释放内存, 让系统保持整体可用。	<p>可以通过以下命令查看参数取值:</p> <pre>cat /proc/sys/vm/panic_on_oom或者 sysctl -a grep panic_on_oom</pre> <ul style="list-style-type: none">• 值为0: 内存不足时, 触发OOM killer。• 值为1: 内存不足时, 根据具体情况可能发生kernel panic, 也可能触发OOM killer。• 值为2: 内存不足时, 强制触发系统panic, 导致系统重启。 <p>说明 HCE中参数默认值为1。</p>	<p>例如将参数设置为0, 可用以下两种方式:</p> <ul style="list-style-type: none">• 临时配置, 立即生效, 但重启后恢复成默认值。 sysctl -w vm.panic_on_oom=0• 持久化配置, 系统重启仍生效。执行命令 vim /etc/sysctl.conf, 在该文件中添加一行 vm.panic_on_oom =0, 再执行命令sysctl -p或重启系统后生效。

参数名称	参数说明	取值	修改方式
oom_kill_allocating_task	<p>当系统选择触发OOM killer，试图结束某些进程时，oom_kill_allocating_task参数会控制选择哪些进程，有以下两种选择：</p> <ul style="list-style-type: none"> • 触发OOM的进程。 • oom_score得分最高的进程。 	<p>可以通过以下命令查看参数取值：</p> <pre>cat /proc/sys/vm/oom_kill_allocating_task或者sysctl -a grep oom_kill_allocating_task</pre> <ul style="list-style-type: none"> • 值为0：选择oom_score得分最高的进程。 • 值为非0：选择触发OOM的进程。 <p>说明 HCE中参数默认值为0。</p>	<p>例如将该参数设置成1，可用以下两种方式：</p> <ul style="list-style-type: none"> • 临时配置，立即生效，但重启后恢复成默认值。 sysctl -w vm.oom_kill_allocating_task=1 • 持久化配置，系统重启仍生效。执行命令vim /etc/sysctl.conf，在该文件中添加一行vm.oom_kill_allocating_task=1，再执行命令sysctl -p或重启系统后生效。
oom_score	<p>指进程的得分，主要有两部分组成：</p> <ul style="list-style-type: none"> • 系统打分，主要是根据该进程的内存使用情况由系统自动计算。 • 用户打分，也就是oom_score_adj，可以自定义。 	<p>可以通过调整oom_score_adj的值进而调整一个进程最终的得分。通过以下命令查看参数取值：</p> <pre>cat /proc/进程id/oom_score_adj</pre> <ul style="list-style-type: none"> • 值为0：不调整oom_score。 • 值为负值：在实际打分值上减去一个折扣。 • 值为正值：增加该进程的oom_score。 <p>说明 oom_score_adj的取值范围是-1000~1000。 若设定成OOM_SCORE_ADJ_MIN或-1000，则表示禁止OOM killer结束该进程。</p>	<p>例如将进程id为2939的进程oom_score_adj参数值设置为1000，可用以下命令：</p> <pre>echo 1000 > /proc/2939/oom_score_adj</pre>

参数名称	参数说明	取值	修改方式
oom_dump_tasks	<p>oom_dump_tasks参数控制OOM发生时是否记录系统的进程信息和OOM killer信息。</p> <p>例如dump系统中所有的用户空间进程关于内存方面的一些信息，包括：进程标识信息、该进程使用的内存信息、该进程的页表信息等，这些信息有助于了解出现OOM的原因。</p>	<p>可以通过以下命令查看参数取值：</p> <pre>cat /proc/sys/vm/oom_dump_tasks或者sysctl -a grep oom_dump_tasks</pre> <ul style="list-style-type: none"> 值为0：OOM发生时不会打印相关信息。 值为非0：以下三种情况会调用dump_tasks打印系统中所有task的内存状况。 <ul style="list-style-type: none"> 由于OOM导致kernel panic。 没有找到需要结束的进程。 找到进程并将其结束的时候。 <p>说明 HCE中参数默认值为1。</p>	<p>例如将该参数设置成0，可用以下两种方式：</p> <ul style="list-style-type: none"> 临时配置，立即生效，但重启后恢复成默认值。 sysctl -w vm.oom_dump_tasks=0 持久化配置，系统重启仍生效。执行命令vim /etc/sysctl.conf，在该文件中添加一行vm.oom_dump_tasks=0，再执行命令sysctl -p或重启系统后生效。

触发 OOM killer 示例

- 您可以参考表9-1设置HCE系统参数，示例配置如下：

```
[root@localhost ~]# cat /proc/sys/vm/panic_on_oom
0
[root@localhost ~]# cat /proc/sys/vm/oom_kill_allocating_task
0
[root@localhost ~]# cat /proc/sys/vm/oom_dump_tasks
1
```

- panic_on_oom=0，表示在发生系统OOM的时候触发OOM killer。
- oom_kill_allocating_task=0，表示触发OOM killer的时候优先选择结束得分高的进程。
- oom_dump_tasks=1，表示系统发生OOM的时候记录系统的进程信息和OOM killer信息。

- 启动测试进程。

在系统中同时启动三个相同的测试进程（test、test1、test2），不断申请新的内存，并将test1的oom_score_adj设置成最大的1000，表示OOM killer优先结束该进程，直至内存耗尽触发系统OOM。

```
[root@localhost ~]# ps -ef | grep test
root 2938 2783 0 19:08 pts/2 00:00:00 ./test
root 2939 2822 0 19:08 pts/3 00:00:00 ./test1
root 2940 2918 0 19:08 pts/5 00:00:00 ./test2
[root@localhost ~]# echo 1000 > /proc/2939/oom_score_adj
[root@localhost ~]# cat /proc/2939/oom_score_adj
1000
```

3. 查看OOM信息。

经过一段时间后系统发生OOM并触发OOM killer，同时在/var/log/messages中打印系统所有进程的内存等信息并结束了test1进程：

```
[root@localhost ~]# ps -ef | grep test
root      2938      2783    0 19:08 pts/2    00:00:06 ./test
root      2940      2918    0 19:08 pts/5    00:00:06 ./test2
root      3383      2861    0 19:41 pts/4    00:00:00 grep --color=auto test
[root@localhost ~]#

Mar 4 19:15:13 localhost kernel: [5865.355934] kernel-fault(0x2) notification starting on CPU 16
Mar 4 19:15:14 localhost kernel: [5865.355938] kernel-fault(0x2) notification finished on CPU 16
Mar 4 19:15:14 localhost kernel: [5865.355944] top invoked oom: gfp_mask=0x0cc0(GFP_KERNEL), order=0
Mar 4 19:15:14 localhost kernel: [5865.355950] oom_score_adj: 0
Mar 4 19:15:14 localhost kernel: [5865.355950] Hardware name: OpenStack Foundation OpenStack Nova, BIOS rel-1.12.1-0-ga5cab58-20230319_150422-sxrtoscli0000 04/01/2014
Mar 4 19:15:14 localhost kernel: [5865.355953] Call Trace:
Mar 4 19:15:14 localhost kernel: [5865.355967] dump_stack+0x57/0x6e
Mar 4 19:15:14 localhost kernel: [5865.355970] oom_show_debug_info.part.0+0x4a/0x131
Mar 4 19:15:14 localhost kernel: [5865.355972] ? printk+0x9/0x23
Mar 4 19:15:14 localhost kernel: [5865.355975] out_of_memory.cold+0x68/0x8b
Mar 4 19:15:14 localhost kernel: [5865.355978] _alloc_pages+0x69/0x1250
Mar 4 19:15:14 localhost kernel: [5865.355980] pagecache_get_page+0x1a3/0x560
Mar 4 19:15:14 localhost kernel: [5865.355982] filemap_fault+0x1a2/0x790
Mar 4 19:15:14 localhost kernel: [5865.356000] ext4_filemap_fault+0x2d/0x50 [ext4]
Mar 4 19:15:14 localhost kernel: [5865.356070] do_fault+0x37/0x1e0
Mar 4 19:15:14 localhost kernel: [5865.356071] do_read_fault+0x11/0x10
Mar 4 19:15:14 localhost kernel: [5865.356072] do_fault+0x9/0x100
Mar 4 19:15:14 localhost kernel: [5865.356074] _handle_mm_fault+0x448/0x7b0
Mar 4 19:15:14 localhost kernel: [5865.356076] handle_mm_fault+0x69/0x210
Mar 4 19:15:14 localhost kernel: [5865.356079] exc_page_fault+0x1db/0x5a0
Mar 4 19:15:14 localhost kernel: [5865.356081] ? asm_exc_page_fault+0x20/0x30
Mar 4 19:15:14 localhost kernel: [5865.356082] asm_exc_page_fault+0x19/0x30
Mar 4 19:15:14 localhost kernel: [5865.356088] RIP: 0033:0x7f4b7b88b360
Mar 4 19:15:14 localhost kernel: [5865.356091] Code: Unable to access opcode bytes at RIP 0x7f4b7b88b360.
Mar 4 19:15:14 localhost kernel: [5865.356094] RSP: 002b:00007ffefed523a0 EFLAGS: 00010202
Mar 4 19:15:14 localhost kernel: [5865.356098] RAX: 00007f4b7b8e6888 RBX: 0000000000000000 RCX: 0000000000000000
Mar 4 19:15:14 localhost kernel: [5865.356099] RDX: 0000000000010001 RSI: 00007f4b7b8e6888 RDI: 0000000000000000
Mar 4 19:15:14 localhost kernel: [5865.356097] RBP: 00007f4b7b897800 R08: 00007ffefed523d0 R09: 00007ffefed523e0
Mar 4 19:15:14 localhost kernel: [5865.356098] R10: 0000000000000240 R11: 0000000000000240 R12: 00007ffefed52450
Mar 4 19:15:14 localhost kernel: [5865.356099] R13: 0000000000000001 R14: 0000000000000000 R15: 00005cfd5ed1d1a1
Mar 4 19:15:14 localhost kernel: [5865.360241] error: MEM_PFN(All)Memory is NULL
Mar 4 19:15:14 localhost kernel: [5865.360270] #PID# info:
Mar 4 19:15:14 localhost kernel: [5865.360280] slabinfo - version: 2.1
Mar 4 19:15:14 localhost kernel: [5865.369240] # name          <active_objs> <num_objs> <objsize> <objperslab> <pagesperslab> : tunables <limit> <batchcount> <sharedfactor>
Mar 4 19:15:14 localhost kernel: [5865.369241] nf_contrack_expect 0 0 248 33 2 : tunables 0 0 0 : slabdata 0 0 0
Mar 4 19:15:14 localhost kernel: [5865.369244] nf_contrack 561 561 320 51 4 : tunables 0 0 0 : slabdata 11 11 0
Mar 4 19:15:14 localhost kernel: [5865.369253] rpc_inode_cache 46 46 704 46 0 : tunables 0 0 0 : slabdata 1 1 0

Mar 4 19:15:15 localhost kernel: [5865.370050] 2938 0 2938 2938 2 80932 297 0 bash
Mar 4 19:15:15 localhost kernel: [5865.370051] 2938 0 2938 656135 285530 5267456 365022 0 test
Mar 4 19:15:15 localhost kernel: [5865.370051] 2938 0 2938 60320 291508 262056 24284 1000 test1
Mar 4 19:15:15 localhost kernel: [5865.370061] 2940 0 2940 40025 20480 4878056 309516 0 test2
Mar 4 19:15:15 localhost kernel: [5865.370062] 2941 0 2941 4555 40 65536 293 0 sshd
Mar 4 19:15:15 localhost kernel: [5865.370063] 2942 0 2942 4680 141 65536 238 0 sshd
Mar 4 19:15:15 localhost kernel: [5865.370065] 2946 0 2946 4577 44 73728 291 0 sshd
Mar 4 19:15:15 localhost kernel: [5865.370066] 2950 0 2950 4577 43 73728 293 0 sshd
Mar 4 19:15:15 localhost kernel: [5865.370068] 2951 0 2951 2827 0 57344 135 0 sftp-server
Mar 4 19:15:15 localhost kernel: [5865.370071] 2951 0 2951 3926 83 69032 324 0 bash
Mar 4 19:15:15 localhost kernel: [5865.370072] 3010 0 3010 6435 139 73728 82 0 top
Mar 4 19:15:15 localhost kernel: [5865.370072] oom-kill:CONSTRAINT_NONE:node=0(full),dps=0r; mem_all=0; global_oom_task_memcg/system.slice/ssh.service;task=test1,pid=2939,uid=0
Mar 4 19:15:15 localhost kernel: [5865.370080] Out of memory: Killed process 2939 (test1) total-vm:252212kB, anon-rss:116072kB, file-rss:4kB, shmem-rss:0kB, UID:0 pgtables:208kB oom_score_adj:1000
Mar 4 19:20:00 localhost systemd[1]: Starting... accounting tool...
Mar 4 19:20:03 localhost systemd[1]: sssstat-collect.service: deactivated successfully.
```

OOM 可能的原因

- **cgroup内存不足**
使用的内存超出了cgroup中memory.limit_in_bytes配置的大小，如下示例演示memory.limit_in_bytes配置为80M，使用memhog模拟分配100M，触发OOM，/var/log/messages部分日志如下，可以从日志中看到memhog所在进程（PID: 2021820）使用了81920kB内存，超出了限制，触发了OOM：
warning|kernel[-][2919920.414131] memhog invoked oom-killer: gfp_mask=0xccc0(GFP_KERNEL), order=0, oom_score_adj=0
info|kernel[-][2919920.414220] memory: usage 81920kB, limit 81920kB, failcnt 30
err|kernel[-][2919920.414272] Memory cgroup out of memory: Killed process 2021820 (memhog) total-vm:105048kB, anon-rss:81884kB, file-rss:1544kB, shmem-rss:0kB, UID:0 pgtables:208kB oom_score_adj:0
- **父cgroup内存不足**
在子cgroup中内存仍然足够，但是父cgroup的内存不足，超过了内存限制，如下示例演示父cgroup memory.limit_in_bytes配置为80M，两个子cgroup memory.limit_in_bytes均配置为50M，在两个子cgroup中使用程序循环分配内存，触发OOM，/var/log/messages部分日志如下：
warning|kernel[-][2925796.529231] main invoked oom-killer: gfp_mask=0xccc0(GFP_KERNEL), order=0, oom_score_adj=0
info|kernel[-][2925796.529315] memory: usage 81920kB, limit 81920kB, failcnt 199
err|kernel[-][2925796.529366] Memory cgroup out of memory: Killed process 3238866 (main) total-vm:46792kB, anon-rss:44148kB, file-rss:1264kB, shmem-rss:0kB, UID:0 pgtables:124kB oom_score_adj:0
- **系统全局内存不足**
一方面由于OS的空闲内存不足，有程序一直在申请内存，另一方面也无法通过内存回收机制解决内存不足的问题，因此触发了OOM，如下示例演示OS中使用程序循环分配内存，触发OOM，/var/log/messages部分日志如下，可以从日志中看到内存节点Node 0的空闲内存（free）已经低于了内存最低水位线（low），触发了OOM：

```
kernel: [ 1475.869152] main invoked oom: gfp_mask=0x100dca(GFP_HIGHUSER_MOVABLE|
__GFP_ZERO), order=0
kernel: [ 1477.959960] Node 0 DMA32 free:22324kB min:44676kB low:55844kB high:67012kB
reserved_highatomic:0KB active_anon:174212kB inactive_anon:1539340kB active_file:0kB
inactive_file:64kB unevictable:0kB writepending:0kB present:2080636kB managed:1840628kB
mlocked:0kB pagetables:7536kB bounce:0kB free_pcp:0kB local_pcp:0kB free_cma:0kB
kernel: [ 1477.960064] oom-
kill:constraint=CONSTRAINT_NONE,nodemask=(null),cpuset=/,mems_allowed=0,global_oom,task_mem
cg=/system.slice/ssh.service,task=main,pid=1822,uid=0
kernel: [ 1477.960084] Out of memory: Killed process 1822 (main) total-vm:742748kB, anon-
rss:397884kB, file-rss:4kB, shmem-rss:0kB, UID:0 pgtables:1492kB oom_score_adj:1000
```

- 内存节点（Node）的内存不足

在NUMA存储模式下，OS会存在多个内存节点，如果程序制定使用特定节点的内存，可能在OS内存充足的情况下触发OOM，如下示例演示在两个内存节点的情况下，使用程序循环在Node 1分配内存，导致Node 1内存不足，但是OS内存足够，/var/log/messages部分日志如下：

```
kernel: [ 465.863160] main invoked oom: gfp_mask=0x100dca(GFP_HIGHUSER_MOVABLE|
__GFP_ZERO), order=0
kernel: [ 465.878286] active_anon:218 inactive_anon:202527 isolated_anon:0#012 active_file:5979
inactive_file:5231 isolated_file:0#012 unevictable:0 dirty:0 writeback:0#012 slab_reclaimable:6164
slab_unreclaimable:9671#012 mapped:4663 shmem:2556 pagetables:846 bounce:0#012 free:226231
free_pcp:36 free_cma:0
kernel: [ 465.878292] Node 1 DMA32 free:34068kB min:32016kB low:40020kB high:48024kB
reserved_highatomic:0KB active_anon:188kB inactive_anon:778076kB active_file:20kB
inactive_file:40kB unevictable:0kB writepending:0kB present:1048444kB managed:866920kB
mlocked:0kB pagetables:2752kB bounce:0kB free_pcp:144kB local_pcp:0kB free_cma:0kB
kernel: [ 933.264779] oom-
kill:constraint=CONSTRAINT_MEMORY_POLICY,nodemask=1,cpuset=/,mems_allowed=0-1,global_oom,
task_memcg=/system.slice/ssh.service,task=main,pid=1733,uid=0
kernel: [ 465.878438] Out of memory: Killed process 1734 (main) total-vm:239028kB, anon-
rss:236300kB, file-rss:200kB, shmem-rss:0kB, UID:0 pgtables:504kB oom_score_adj:1000
```

- 其他可能原因

OS在内存分配的过程中，如果伙伴系统的内存不足，则系统会通过OOM Killer释放内存，并将内存提供至伙伴系统。

OOM 问题解决方法

- 从业务进程排查，确认是否有内存泄漏，导致OOM。
- 排查cgroup limit_in_bytes配置是否与业务内存规划匹配，如需要调整，可以手动执行以下命令修改配置参数：
echo <value> > /sys/fs/cgroup/memory/<cgroup_name>/memory.limit_in_bytes
- 如果确认业务需要比较多的内存，建议升级弹性云服务器内存规格。

10 IPVS 报错问题说明

问题背景

IPVS (IP Virtual Server) 指IP虚拟服务器，用于负载均衡、网络转发等目的。用户在系统上配置了IPVS虚拟服务器，但未配置真实服务器的情况下，会在VNC上出现错误日志。

问题现象

配置了IPVS虚拟服务器，但未配置真实服务器时，当网络请求发往该虚拟服务器地址后，通过华为云VNC登录的控制台上可以看到类似如下的错误日志。

```
[32264.645949][T268365] IPVS: wlc: TCP 192.168.1.1:5000 - no destination available
[32265.234919][T268366] IPVS: wlc: TCP 192.168.1.1:5000 - no destination available
[32265.954662][T268367] IPVS: wlc: TCP 192.168.1.1:5000 - no destination available
[32266.557032][T268368] IPVS: wlc: TCP 192.168.1.1:5000 - no destination available
[32267.166530][T268369] IPVS: wlc: TCP 192.168.1.1:5000 - no destination available
[32267.725920][T268370] IPVS: wlc: TCP 192.168.1.1:5000 - no destination available
```

解决方法

1. 安装ipvsadm。
2. 执行ipvsadm -Ln，查询当前虚拟服务器的配置。找到报错的虚拟服务器对应的表项。

图 10-1 未配置真实服务器

```
IP Virtual Server version 1.2.1 (size=4096)
Prot LocalAddress:Port[,Subnet] Scheduler Established(Sec.) Flags
-> RemoteAddress:Port[,:0if] Forward Weight ActiveConn InActConn UtepAddr:vtepport Unid Mac
TCP 192.168.1.1:5000 wlc
```

如上图所示没有对应的真实服务器，则说明配置不完整，会引发错误打印。需要排查对应的业务流程是否正确。

图 10-2 已配置真实服务器

```
IP Virtual Server version 1.2.1 (size=4096)
Prot LocalAddress:Port[,Subnet] Scheduler Established(Sec.) Flags
-> RemoteAddress:Port[,:0if] Forward Weight ActiveConn InActConn UtepAddr:vtepport Unid Mac
TCP 192.168.1.1:5000 wlc
-> 192.168.1.2:5000 Masq 1 0 0
```

如上图是完整的有真实服务器的配置。

3. 如果希望排除该IPVS错误日志在VNC上对用户操作的干扰，可以按如下的方式处理（选择其中一种即可）

- 关闭业务发送的网络请求，具体操作需要用户根据自身业务情况来处理。
- 执行以下命令调整内核printk打印等级。

```
echo 3 4 1 7 > /proc/sys/kernel/printk
```

说明

如果临时修改系统配置，建议用户选择适当的时机恢复系统配置。

- 用华为云的CloudShell方式登录云服务器进行操作。

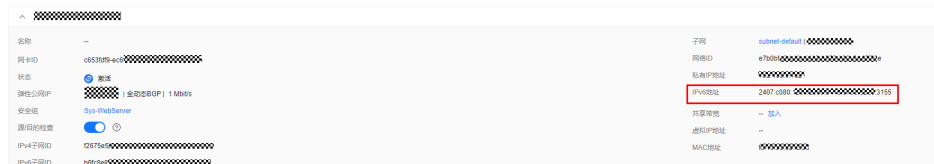
11 ECS 开启 IPv6 后，HCE OS 系统内无法获取到 IPv6 地址

问题背景

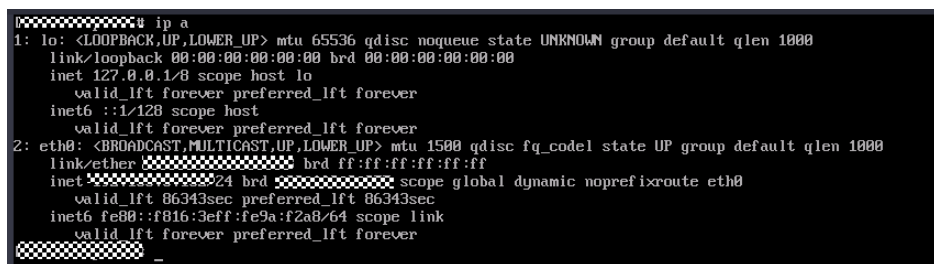
在弹性云服务器ECS控制台上开启云服务器网卡的IPv6功能后，由于未在操作系统内部正确配置IPv6，导致HCE OS系统内无法获取到IPv6地址。

问题现象

在弹性云服务器ECS控制台上已开启IPv6功能，在详情界面已显示IPv6地址。



但是进入操作系统内部，无法获取到IPv6地址。



解决方法

1. 手动配置dhcp自动获取ipv6地址，如下图在对应网卡配置文件（/etc/sysconfig/network-scripts/ifcfg-ethx）中添加以下参数。

IPV6INIT="yes"

DHCPV6C="yes"

```
[root@~]# cat /etc/sysconfig/network-scripts/ifcfg-eth0
DEVICE="eth0"
BOOTPROTO="dhcp"
ONBOOT="yes"
TYPE="Ethernet"
PERSISTENT_DHCLIENT="yes"
IPV6INIT="yes"
DHCPV6C="yes"
[root@~]#
```

2. 执行以下命令重启NetworkManager服务即可获取到IPv6地址。

systemctl restart NetworkManager

```
[root@~]# systemctl restart NetworkManager
[root@~]# ip a
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN group default qlen 1000
    link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
    inet 127.0.0.1/8 scope host lo
        valid_lft forever preferred_lft forever
    inet6 ::1/128 scope host
        valid_lft forever preferred_lft forever
2: eth0: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc fq_codel state UP group default qlen 1000
    link/ether 00:00:00:00:00:00 brd ff:ff:ff:ff:ff:ff
    inet 192.168.1.100/24 brd 192.168.1.255 scope global dynamic noprefixroute eth0
        valid_lft 86394sec preferred_lft 86394sec
    inet6 2407::3155/128 scope global dynamic noprefixroute
        valid_lft 7195sec preferred_lft 7195sec
    inet6 fe80::f816:3eff:fe9a:f2a8/64 scope link noprefixroute
        valid_lft forever preferred_lft forever
[root@~]#
```